

EFFICIENT PROCESS FOR HANDOVER BETWEEN SUBNET MANAGERS

5 The present invention is related to the subject
matter of the following commonly assigned, co-pending
United States Patent Applications filed concurrently
herewith: Serial No. _____ (Docket No. AUS9-2000-0620)
entitled "Method and System for Informing An Operating
System In A System Area Network When A New Device Is
10 Connected"; Serial No. _____ (Docket No. AUS9-2000-
0622) entitled "Method and System For Scalably Selecting
Unique Transaction Identifiers"; Serial No. _____
(Docket No. AUS9-2000-0623) entitled "Method And System
For Reliably Defining and Determining Timeout Values In
Unreliable Datagrams"; and Serial No. _____ (Docket No.
AUS9-2000-0624) entitled "Method and System For Choosing
A Queue Protection Key That is Tamper-proof From An
Application". The content of the above-referenced
applications is incorporated herein by reference.

BACKGROUND OF THE INVENTION**1. Technical Field:**

25 The present invention relates in general to computer
networks and, in particular, to merging of independent
computer networks. Still more particularly, the present
invention relates to a method and system for providing an
30 efficient handover of control between Subnet managers of
separate subnets, which are being merged into a single
subnet.

35 *Cancel*

2. Description of the Related Art:

The use of I/O interconnects to connect components of a distributed computer system is known in the art. Traditionally, in such systems, individual components are interconnected via a parallel bus, such as a PCIX bus. The parallel bus has a relatively small number of plug-in ports for connecting the components. The number of plug-in-ports is set (i.e., the number cannot be increased). At maximum loading, a PCIX bus transmits data at about 1 Gbyte/second.

The introduction of high performance adapters (e.g., SCSI adapters), Internet-based networks, and other high performance network components has resulted in increased demand for bandwidth, faster network connections, distributed processing functionality, and scaling with processor performance. These and other demands are quickly outpacing the current parallel bus technology and are making the limitations of parallel buses even more visible. PCIX bus, for example, is not scalable, i.e., the length of the bus and number of slots available at a given frequency cannot be expanded to meet the needs for more components, and the limitation hinders further development of fast, efficient distributed networks, such as system area networks. New switched network topologies and systems capable of being easily expanded are required to keep up with the increasing demands, while allowing the network processes on the expanding network to be dynamically completed. i.e., without manual input.

The present invention recognizes the need for

faster, more efficient computer interconnects offering the features demanded by the developments of technology. More specifically, the present invention recognizes the need for providing a mechanism within a network such as a System Area Network (SAN) consisting of multiple subnets, that provides efficient, dynamic combining of two or more subnets into a single network.

2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
26

SUMMARY OF THE INVENTION

A method and system for efficiently combining (or merging) subnets having individual master subnet managers into a single network with one master subnet manager is disclosed. The invention is thus applicable to a distributed computing system, such as a system area network, having end nodes, switches, and routers, and links interconnecting these components. The switches and routers interconnect the end nodes and route packets, i.e., sub-components of messages being transmitted, from a source end node to a target end node. The target end node then reassembles the packets into the message.

During discovery and configuration of a subnet, a subnet manager creates a Subnet Management Database (SMDB) representative of the subnet components being managed. Each subnet manager contains an independent SMDB. When two or more subnets are merged (i.e., linked/connected together) to form a single network, a single one of the subnet managers is selected as the master subnet manager and all the subnets SMDBs must be merged. The other subnet managers are relegated to standby status. In a preferred embodiment, a SMDB record labeling mechanism is utilized to differentiate among components from the different subnets that may have the same parameter values, such as protection keys (P_keys).

In one embodiment of the invention, when there are two or more separate SMDBs, each maintained by separate master subnet managers on a separate subnet, when those separate subnets are linked together requiring they be merged into

one larger subnet, means are provided for efficient merging the SMDBs. In this manner the separate subnets become a single subnet, with a single master subnet manager and with a single SMDB.

5

All objects, features, and advantages of the present invention will become apparent in the following detailed written description.

[illegible]

BRIEF DESCRIPTION OF THE DRAWINGS

The novel features believed characteristic of the invention are set forth in the appended claims. The invention itself however, as well as a preferred mode of use, further objects and advantages thereof, will best be understood by reference to the following detailed description of an illustrative embodiment when read in conjunction with the accompanying drawings, wherein:

Figure 1A depicts a system area network (SAN) in which the present invention is preferably implemented;

Figure 1B is a diagram illustrating the software aspects of SAN management model in accordance with the present invention;

Figure 2 is a diagram illustrating software aspects of an exemplary host processor end node for the SAN of **Figure 1** in accordance with the present invention;

Figure 3 is a diagram of an exemplary host channel adapter of the SAN of **Figure 1** in accordance with the present invention;

Figure 4 is a diagram illustrating processing of work requests with queues in accordance with a preferred embodiment of the present invention;

Figure 5 is an illustration of a data packet in accordance with a preferred embodiment of the present invention;

Figure 6 is a diagram illustrating a communication over a portion of a SAN fabric;

Figure 7 is a diagram illustrating packet transfers in accordance with the invention;

Figure 8 is a diagram illustrating two independent networks, which are linked together to form a single network in accordance with a preferred embodiment of the invention;

Figure 9A is a flow chart showing the process of time stamping SMDB entries in accordance with a preferred embodiment of the invention; and

Figure 9B is a flow chart showing the process absorbing an SMDB by a master subnet manager in accordance a preferred embodiment of the invention.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

The present invention is directed to a method for efficient handover of control to a single master subnet manager when two or more subnets having subnet managers are merged. The invention is applicable to a distributed computer system, such as a system area network (SAN). The invention is implemented in a manner that allows the merged subnet overlaps in parameters, such as P_Keys, to be resolved.

In order to appreciate the environment within which the invention is preferably practiced, a description of a SAN configured with routers, switches, and end nodes, etc. is provided below. Presentation of the environment and particular functional aspects of the environment which enable the invention to be practiced are provided with reference to **Figures 1-5**. Section headings have been provided to distinguish the hardware and software architecture of the SAN. However, those skilled in the art understand that the descriptions of either architecture necessarily includes references to both components.

SAN HARDWARE ARCHITECTURE

With reference now to the figures and in particular with reference to **Figure 1**, there is illustrated an exemplary embodiment of a distributed computer system. Distributed computer system **100** represented in **Figure 1** is provided merely for illustrative purposes, and the embodiments of the present invention described below can

be implemented on computer systems of numerous other types and configurations. For example, computer systems implementing the present invention may range from a small server with one processor and a few input/output (I/O) adapters to very large parallel supercomputer systems with hundreds or thousands of processors and thousands of I/O adapters. Furthermore, the present invention can be implemented in an infrastructure of remote computer systems connected by an Internet or intranet.

As shown in **Figure 1**, distributed computer system **100** includes a system area network (SAN) **113**, which is a high-bandwidth, low-latency network interconnecting nodes within the distributed computer system. More than one (1) SAN **113** may be included in a distributed computer system **100** and each SAN **113** may comprise multiple sub-networks (subnets).

A node is herein defined to be any component that is attached to one or more links of a network. In the illustrated distributed computer system, nodes include host processors **101**, redundant array of independent disks (RAID) subsystem **103**, I/O adapters **105**, switches **109A-109C**, and router **111**. The nodes illustrated in **Figure 1** are for illustrative purposes only, as SAN **113** can connect any number and any type of independent nodes. Any one of the nodes can function as an end node, which is herein defined to be a device that originates or finally consumes messages or frames in the distributed computer system **100**.

SAN **113** is the communications and management infrastructure supporting both I/O and inter-processor communications (IPC) within distributed computer system **100**. Distributed computer system **100**, illustrated in **Figure 1**, includes a switched communications fabric (i.e., links, switches and routers) allowing many devices to concurrently transfer data with high-bandwidth and low latency in a secure, remotely managed environment. End nodes can communicate over multiple ports and utilize multiple paths through SAN **113**. The availability of multiple ports and paths through SAN **113** can be employed for fault tolerance and increased-bandwidth data transfers.

SAN **113** includes switches **109A-109C** and routers **111**. Switch **109A-109C** connects multiple links together and allows routing of packets from one link to another link within SAN **113** using a small header Destination Local Identifier (DLID) field. Router **111** is capable of routing frames from one link in a first subnet to another link in a second subnet using a large header Destination Globally Unique Identifier (DGUID). Router **111** may be coupled via wide area network (WAN) and/or local area network (LAN) connections to other hosts or other routers.

In SAN **113**, host processor nodes **101** and I/O nodes **106** include at least one Channel Adapter (CA) to interface to SAN **113**. Host processor nodes **101** include central processing units (CPUs) **119** and memory **121**. In one embodiment, each CA is an endpoint that implements

the CA interface in sufficient detail to source or sink packets transmitted on SAN **113**. As illustrated, there are two CA types, Host CA (HCA) **117** and Target CA (TCA) **127**. HCA **117** is used by general purpose computing nodes to access SAN **113**. In one implementation, HCA **117** is implemented in hardware. In the hardware implementation of HCA **117**, HCA hardware offloads much of CPU and I/O adapter communication overhead. The hardware implementation of HCA **117** also permits multiple concurrent communications over a switched network without the traditional overhead associated with communicating protocols. Use of HCAs **117** in SAN **113** also provides the I/O and IPC consumers of distributed computer system **100** with zero processor-copy data transfers without involving the operating system kernel process. HCA **117** and other hardware of SAN **113** provide reliable, fault tolerant communications.

The I/O chassis **106** includes I/O adapter backplane and multiple I/O adapter nodes **105** that contain adapter cards. Exemplary adapter cards illustrated in **Figure 1** include SCSI adapter card **123A**, adapter card **123B** to fiber channel hub and FC-AL devices, Ethernet adapter card **123C**, graphics adapter card **123D**, and video adapter card **123E**. Any known type of adapter card can be implemented. The I/O chassis **106** also includes switch **109B** in the I/O adapter backplane to couple adapter cards **123A-123E** to SAN **113**.

RAID subsystem **103** includes a microprocessor **125**, memory **126**, a Target Channel Adapter (TCA) **127**, and

multiple redundant and/or striped storage disks **129**.

In the illustrated SAN **113**, each link **115** is a full duplex channel between any two network elements, such as end nodes, switches **109A-109C**, or routers **111**. Suitable links **115** may include, but are not limited to, copper cables, optical cables, and printed circuit copper traces on backplanes and printed circuit boards. The combination of links **115** and switches **109A-109C**, etc. operate to provide point-to-point communication between nodes of SAN **113**.

SAN SOFTWARE ARCHITECTURE

Software Components

Software and hardware aspects of an exemplary host processor node **101** are generally illustrated in **Figure 2**. Host processor node **101** includes one or more processors that execute a set of consumer processes **201**. Host processor node **101** includes HCA **117** with ports **205**. Each port **205** connects to a link **115** of SAN **113**. Ports **205** can connect to one SAN subnet or multiple SAN subnets. Utilizing message and data services **203**, consumer processes **201** transfer messages to SAN **113** via verbs interface **207**. Verbs interface **207** is generally implemented with an operating-system specific programming interface.

A software model of HCA **117** is illustrated in **Figure 3**. HCA **117** includes a set of queue pairs (QPs) **301**, which transfer messages across ports **205** to the subnet.

A single HCA **117** may support thousands of QPs **301**. By contrast, TCA **127** in an I/O adapter typically supports a much smaller number of QPs **301**. Also illustrated are subnet management administration (SMA) **209**, management packets **211** and a number of virtual lanes **213**, which connect transport layer with ports **205**.

Turning now to **Figure 5**, there is illustrated a software management model for nodes on SAN **113**. SAN architecture management facilities provides a Subnet Manager (SM) **303A**, a Subnet Administration (SA) **303B**, and an infrastructure that supports a number of general management services. The management infrastructure includes a Subnet Management Agent (SMA) **307** operating in each node and defines a general service interface that allows additional general services agents. Also, SAN architecture defines a common management datagram (MAD) message structure for communicating between managers and management agents.

SM **303A** is responsible for initializing, configuring and managing switches, routers, and channel adapters. The SM can be implemented within other devices, such as a channel adapter or a switch. One SM **303A** of SAN is dedicated as a master SM and is responsible for: discovering the subnet topology; configuring each channel adapter port with a range of Local Identification (LID) numbers, Global Identification (GID) number, subnet prefix, and Partition Keys (P_Keys); configuring each switch with a LID, the subnet prefix, and with its forwarding database; and maintaining the end node and

service databases for the subnet to provide a Global Unique Identification (GUID) number to LID/GID resolution service as well as a services directory. Thus, management of SAN **113** and SAN components, such as HCAs **117**, TCAs (or end nodes) **127**, switches **109**, and routers **111** are completed utilizing Subnet Management (SM) **303A** and Subnet Administration (SA) **303B**. SMPs are used to discover, initialize, configure, and maintain SAN components through management agents **307** of end nodes **305**. SAN SA packets are used by SAN components to query and update subnet management data. Control of some aspects of the subnet management are provided via a user management console **311** in host-based end node **309**.

MESSAGE TRANSFER PROCESS

SAN **113** provides the high-bandwidth and scalability required for I/O and also supports the extremely low latency and low CPU overhead required for Interprocessor Communications (IPC). User processes can bypass the operating system (OS) kernel process and directly access network communication hardware, such as HCAs **117**, which enable efficient message passing protocols. SAN **113** is suited to current computing models and is a building block for new forms of I/O and computer cluster communication. SAN **113** allows I/O adapter nodes **105** to communicate among themselves or communicate with any or all of the processor nodes **101** in the distributed computer system. With an I/O adapter attached to SAN **113**, the resulting I/O adapter node **105** has substantially the same communication capability as any processor node

101 in the distributed computer system.

For reliable service types of messages, end nodes, such as host processor nodes 101 and I/O adapter nodes 105, generate request packets and receive acknowledgment packets. Switches 109A-109C and routers 111 pass packets along from the source to the target (or destination). Except for the variant CRC trailer field, which is updated at each transfer stage in the network, switches 109A-109C pass the packets along unmodified. Routers 111 update the variant CRC trailer field and modify other fields in the header as the packet is routed.

In SAN 113, the hardware provides a message passing mechanism that can be used for Input/Output (I/O) devices and Interprocess Communications (IPC) between general computing nodes. Consumers (i.e., processing devices connected to end nodes) access SAN 113 message passing hardware by posting send/receive messages to send/receive work queues (WQ), respectively, on a SAN Channel Adapter (CA).

A message is herein defined to be an application-defined unit of data exchange, which is a primitive unit of communication between cooperating processes. A packet (or frame) is herein defined to be one unit of data encapsulated by networking protocol headers (and trailer). The headers generally provide control and routing information for directing the packet (or frame) through SAN 113. The trailer generally contains control and cyclic redundancy check (CRC) data for ensuring that frames are not delivered with corrupted content.

Consumers use SAN verbs to access HCA functions. The software that interprets verbs and directly accesses the CA is known as the Channel Interface (CI) **219**. Send/Receive work queues (WQ) are assigned to a consumer as a Queue Pair (QP). Messages may be sent over five different transport types, Reliable Connected (RC), Reliable Datagram (RD), Unreliable Connected (UC), Unreliable Datagram (UD), and Raw Datagram (RawD). Consumers retrieve the results of these messages from a Completion Queue (CQ) through SAN send and receive work completions (WC). The source CA takes care of segmenting outbound messages and sending them to the destination. The destination or target CA takes care of reassembling inbound messages and placing them in the memory space designated by the destination's consumer. These features are illustrated in the figures below.

Referring now to **Figure 4**, there is illustrated a block diagram of work and completion queue processing. Each QP **301** provides an input to a Send Work Queue (SWQ) **407** and a Receive Work Queue (RWQ) **409**. SWQ **407** sends channel and memory semantic messages, and RWQ **409** receives channel semantic messages. A consumer calls a verb (within verbs interface **207**) to place Work Requests (WRs) into a WQ. A Send WR **403** is a channel semantic operation to push a set of local data segments **417** to the data segments referenced by a remote node's Receive WQE **405**. Each of the Send WR's data segments **417** contains a virtually contiguous memory region. The virtual addresses used to reference the local data segments **417** are in the address context of the process that created

the local QP 301.

As shown in **Figure 4**, WRs 403 that have been placed onto a WQ by consumer process 401 are referred to as work queue elements (WQEs) 405. WQEs 405 are executed by hardware 415 in HCA 117. SWQ 407 contains WQEs 405 that describe data to be transmitted on the SAN fabric. RWQ 409 contains WQEs 405 that describe where to place incoming channel semantic data received from SAN 113.

In one embodiment, Receive Work Queues 409 only support one type of WQE 405, which is referred to as a receive WQE. The receive WQE provides a channel semantic operation describing a local memory space into which incoming send messages are written. The receive WQE includes a scatter list describing several virtually contiguous memory spaces. An incoming send message is written to these memory spaces. The virtual addresses are in the address contexts of the process that created the local QP 301.

Verbs interface 207 also provide a mechanism for retrieving completed work from completion queue 411. Completion queue 411 contains Completion Queue Elements (CQEs) 413 which 413 contain information about previously completed WQEs 405. Completion queues 413 are employed to create a single point of completion notification for multiple QPs 301. CQE 413 contains sufficient information to determine the QP 301 and specific WQE 405 that completed. A completion queue context (not shown) is a block of information that contains pointers to,

length, and other information needed to manage individual completion queues 411.

REMOTE OPERATION FUNCTIONALITY

5

SAN 113, with its interlinked arrangement of components and sub-components, provides a method for completing remote operations, by which processor nodes may directly control processes in I/O nodes. Remote operation also permits the network to manage itself. A remote direct memory access (RDMA) Read WR provides a memory semantic operation to read a virtually contiguous memory space on a remote node. A memory space can either be a portion of a memory region or a portion of a memory window. A memory region references a previously registered set of virtually contiguous memory addresses defined by a virtual address and length. A memory window references a set of virtually contiguous memory addresses which have been bound to a previously registered region.

20

The RDMA Read WR writes the data to a virtually contiguous local memory space. Similar to Send WR 403, virtual addresses used by the RDMA Read WQE to reference the local data segments are in the address context of the process that created the local QP 301. The remote virtual addresses are in the address context of the process owning the remote QP targeted by the RDMA Read WQE.

25

30

RDMA Write WQE provides a memory semantic operation to write a virtually contiguous memory space on a remote node. RDMA Write WQE contains a scatter list of local

virtually contiguous memory spaces and the virtual address of the remote memory space into which the data from the local memory spaces is written.

5 RDMA FetchOp WQE provides a memory semantic operation to perform an atomic operation on a remote word. RDMA FetchOp WQE is a combined RDMA Read, Modify, and Write operation. RDMA FetchOp WQE can support several read-modify-write operations, such as "Compare and Swap if Equal."

10 A Bind (unbind) remote access key (R_Key) WQE provides a command to the HCA hardware to modify a memory window by associating the memory window to a memory region. A second command to destroy a memory window by disassociating the memory window to a memory region is also provided. The R_Key is part of each RDMA access and is used to validate that the remote process has permitted access to the buffer.

20 EFFICIENT PROCESS FOR HANDOVER BETWEEN SUBNET MANAGERS

The present invention makes use of the features described in the subject matter of Patent Application Serial No. _____ (Attorney Docket No AUS9-2000-0625US1) "ASSOCIATION OF END-TO-END CONTEXT VIA RELIABLE DATAGRAM DOMAINS" filed on _____, the entire content of which is hereby incorporated by reference. The reference application allows Reliable Datagram QPs to be used for communicating across multiple partitions. In SAN 113, QPs that support SAN Service Types are associated to a partition and cannot communicate to QPs that are outside of the partition to which the QP is

associated. The QP is barred from communicating with QPs in another partition even if the node's HCA port, which the QP uses, has access to different partitions. RD QPs, however, can communicate with any given partition the node's HCA has access to, so long as there is an underlying End-End Context (EEC) which is associated with the given partition.

The present invention provides a method for seamlessly merging (i.e., linking or connecting) two independent subnets having separate subnet managers (SMs) into a single SAN **113** with a master SM selected from among the separate SMs. The invention utilizes time stamps and GUID addresses on SMDB entries as an aid for a Subnet Manager to efficiently absorb the SMDB of another subnet manager during the subnet management handover process.

Figure 8 illustrates a SAN **113** comprising two separate systems (subnets), each having its own master subnet managers, which are interlinked. System 1 **801** includes Host1 **803**, Master Subnet Manager SM1 **805**, and devices 1A **809** and 1B **807** all interconnected via subnet **810**. System 2 **811** includes Host2 **813**, Master Subnet Manager SM2 **815**, and devices 2A **819** and 2B **817** connected to subnet **820**. In the illustrative embodiment, both systems utilize the same P_Key assignments. Device A1 **809** has a P_Key assigned with a value of 1, and device 1B **807** has a P_Key assigned with a value of 5. Device 2A **819** has a P_Key assigned with a value of 2, and device 2B **817** has a P_Key assigned with a value of 5. Thus, the

two systems are joined together via a link **821** to from
SAN **113**. The linking of the two systems may be completed
via a switch as described above with reference to **Figure**
1A. Thus, devices **1B 807** and **2B 817** of different systems
are each assigned a P_Key value of 5. In a preferred
embodiment, the invention provides a method for handling
the overlaps which occur in the P_Key assignments when
the systems are merged as described below.

When the two systems of **Figure 8** are initially merged,
one of the two master subnet managers is relegated to a
standby subnet manager, while the other SM becomes the
master subnet manager of the new merged subnet (SAN **113**).
Selection of the particular SM that becomes the master SM
of the merged system may be dependent on priority values
associated with each SM during their initial configuration.
In a priority scheme, the SM with the highest priority is
automatically selected as the master SM. Priority value
may be assigned by a subnet administrator when the subnet
is initially configured. In the event that two SMs have
the same priority value, then, in the preferred embodiment,
the SM with the lowest GUID is selected to be the master
subnet manager. During the handover process, the master
subnet manager of the newly-formed merged network receives
the SMDB of the other subnet manager, reconfigures the
network or network components where necessary and takes
over management of the total merged network.

The method utilized by the present invention to
accomplish the seamless transfer of SMDB is to timestamp
each SMDB with the time at which the last change was made
in that SMDB. With this timestamp, the SM that is
accepting the handover can then determine the correct

action to take in handling each entry as the SMDB(s) being taken over is being absorbed. Referring to **Figure 9A**, the process of time stamping entries in SMDB is illustrated. The process begins at block **901**, and following, a determination is made whether there is a change to a Subnet Management Database at block **903**. If not then the process loops back to **901**. If there was a change **903**, then the entry is time stamped with the time of the last change to the entry in that SMDB, and the process returns to block **901**.

Referring now to **Figure 9B**, the process begins at block **951**. Next, a determination is made at block **953** whether there is a merge of subnets occurring, and if no merge is occurring, the process returns to block **951**. If, however, a merge of subnets is occurring, the master subnet manager of the new subnet being formed by merging two or more other subnets begins to examine all entries of SMDBx with it's own SMDBb at block **955**. As utilized herein, SMDBx refers to the database of subnets whose subnet managers are not selected as the master subnet manager when a merger of the subnets occurs. SMDBb refers to the database of the master subnet manager prior to merger of the databases. returning to **Figure 9B**, a determination is made at block **957** whether the master subnet manager finds the same GUID entry in both SMDBx and SMDBb. If the same GUID entry is found, the master subnet manager keeps the one with the latest time stamp and discards the other at block **959**. If, however, the same GUIDs are not found, the process moves to block **961**.

Next a determination is made at block **961** whether all GUIDS have been examined in SMDBx. If not all the GUIDs have been examined, then the process returns to block **957**; however, if all the GUIDs have been examined, the process moves to block **963**, where a determination is made whether SMDBx has the same P_Key entries for different GUIDS as in SMDBb. If yes, then the master subnet manager will change all occurrences of that P_Key in SMDBx to a new unique P_Key value, which is different than that in either SMDBx or SMDBb at block **965**. The process then moves to block **967**. If the P_Key entries are not the same for SMDBx and SMDBb, then a determination is made at block **967** whether all entries in SMDBx have been examined. If all the entries have not been examined, the process returns to block **963**. If all entries have been examined in SMDBx, then the SMDBx is merged with SMDBb at block **969**. Next, a determination is made at block **971** whether there is another database to be merged with master subnet manager's SMDBb, and if so, the process returns to block **955**. Otherwise, the process ends at block **973**.

There may be other overlapping information in the respective databases that may need to be merged as described for the P_Keys and GUIDs above. These would follow a similar process as described in Figure 9B.

As a final matter, it is important to note that while an illustrative embodiment of the present invention has been, and will continue to be, described in the context of a fully functional data processing system, those skilled in

the art will appreciate that the software aspects of an illustrative embodiment of the present invention are capable of being distributed as a program product in a variety of forms, and that an illustrative embodiment of the present invention applies equally regardless of the particular type of signal bearing media used to actually carry out the distribution. Examples of signal bearing media include recordable type media such as floppy disks, hard disk drives, CD ROMs, and transmission type media such as digital and analogue communication links.

In the above detailed description of the preferred embodiments, reference is made to the accompanying drawings which form a part hereof, and in which is shown, by way of illustration, specific embodiments in which the invention may be practiced. It is to be understood that other embodiments may be utilized and structural or logical changes may be made without departing from the scope of the present invention. For example, although the invention is described with reference to multiple computing nodes accessing a single database, the invention is applicable to all other transactions occurring on the network. The above detailed description, therefore, is not to be taken in a limiting sense, and the scope of the present invention is defined by the appended claims.